

## CLAIMS

- 1 1. Apparatus for providing high-performance, scaleable data processing and  
2 storage services to a client from a plurality of resources, comprising  
3 an access interface module which receives requests for service from the  
4 client and selects a subset of the plurality of resources to provide the requested  
5 service and distribute the workload across the plurality of resources; and  
6 a switch fabric for temporarily connecting the access interface module to  
7 the selected subset of the plurality of resources for providing the service to the  
8 client.
- 1 2. The apparatus of claim 1 wherein the access interface module selects the subset  
2 of the plurality of resources based on the relative demand placed on the subset  
3 of resources.
- 1 3. The apparatus of claim 1 wherein the switch fabric comprises a control switch  
2 fabric for transferring control information and a separate data switch fabric for  
3 transferring data.
- 1 4. The apparatus of claim 3 wherein the control switch fabric is optimized for  
2 transferring control information and the data switch fabric is optimized for  
3 transferring data.
- 1 5. The apparatus of claim 3 wherein the request for service includes control  
2 information and data and wherein the access interface module separates the  
3 control information and the data and transfers the data to the selected subset of  
4 resources over the data switch fabric.

- 1 6. The apparatus of claim 3 wherein the data switch fabric comprises a non-  
2 blocking crossbar switch for data transfer and the control switch fabric comprises  
3 an Ethernet switch for control information transfer.
- 1 7. The apparatus of claim 1 further comprising a resource module connected to the  
2 plurality of resources for generating preallocation information that preallocates  
3 services from the plurality of resources in order to evenly distribute a workload  
4 across the plurality of resources.
- 1 8. The apparatus of claim 7 wherein the switch fabric connects the access interface  
2 module to the resource module so that the resource module can transfer the  
3 preallocation information to the access interface module.
- 1 9. The apparatus of claim 8 wherein the access interface module selects a subset  
2 of the plurality of resources based on the preallocation information.
- 1 10. The apparatus of claim 1 wherein the access interface module comprises a data  
2 memory which temporarily stores information transferred between the access  
3 interface module and the selected subset of the plurality of resources.
- 1 11. The apparatus of claim 1 further comprising a plurality of access interface  
2 modules each access interface module receiving service requests from a plurality  
3 of clients.
- 1 12. A disk-based storage system for providing high-performance, scaleable storage  
2 services to a client from a plurality of disks, comprising  
3 a disk interface module connected to the plurality of disks for controlling  
4 data stored on the plurality of disks;

5 a host interface module which receives requests for storage service from  
6 the client and selects a subset of the plurality of disks to provide the requested  
7 storage and distribute the workload across the plurality of disks; and  
8 a switch fabric for temporarily connecting the host interface module to the  
9 resource module for providing the storage service to the client.

1 13. The storage system of claim 12 wherein the switch fabric comprises a control  
2 switch fabric optimized for transferring control information and a separate data  
3 switch fabric optimized for transferring data.

1 14. The storage system of claim 13 wherein the request for service includes control  
2 information and data and wherein the host interface module separates the control  
3 information and the data and transfers the data to the selected subset of  
4 resources over the data switch fabric.

1 15. The storage system of claim 13 wherein the data switch fabric comprises a non-  
2 blocking crossbar switch for data transfer and the control switch fabric comprises  
3 an Ethernet switch for control information transfer.

1 16. The storage system of claim 12 wherein the disk interface module generates  
2 preallocation information that preallocates physical storage in the plurality of  
3 disks in order to evenly distribute data across the plurality of disks.

1 17. The storage system of claim 16 wherein the physical storage in the plurality of  
2 disks is divided into zones and the disk interface module preallocates selected  
3 zones to frequently-accessed data, wherein the selected zones are selected in  
4 order to decrease disk access time.

- 1 18. The storage system of claim 16 wherein the host interface module logically maps  
2 data items to be stored into allocation units preallocated to the host interface  
3 module by the disk interface modules.
- 1 19. The storage system of claim 12 wherein the host interface module comprises a  
2 first data memory and the resource module comprises a second data memory  
3 and wherein the first and second data memories temporarily store information  
4 transferred between the host interface module and the disk interface module.
- 1 20. The storage system of claim 12 further comprising a plurality of host interface  
2 modules, each host interface module receiving service requests from a plurality  
3 of clients.
- 1 21. The storage system of claim 12 further comprising a plurality of disk interface  
2 modules connected to the plurality of disks.
- 1 22. The storage system of claim 21 wherein the host interface module directs writes  
2 to disk interface modules of data based on the relative demand placed on those  
3 modules without regard to either the contents of associated disks or the location  
4 of other copies of the data.
- 1 23. The storage system of claim 22 further comprising a metadata module that is  
2 connected to the switch fabric and assigns each data object an identifying  
3 handle.
- 1 24. The storage system of claim 23 wherein the host interface module and the  
2 metadata module store metadata information associated with each data object in  
3 the plurality of disks separately from, and independently of, data associated with  
4 each data object.

- 1 25. The storage system of claim 24 wherein the handle contains the location of the  
2 metadata in the plurality of disks.
- 1 26. The storage system of claim 23 wherein the metadata module is dedicated  
2 exclusively to fetching, caching and manipulating file and directory attributes of  
3 each data object thereby allowing data paths to be optimized to maximize data  
4 flow.
- 1 27. The storage system of claim 12 wherein the disk interface module preallocates  
2 storage in the plurality of disks in allocation units.
- 1 28. The storage system of claim 27 wherein the host interface module selects a  
2 subset of the plurality of disks from the preallocated storage.
- 1 29. The storage system of claim 27 wherein the host interface module reads a file  
2 from the plurality of disks and writes the file to disk in a single allocation unit.
- 1 30. The storage system of claim 27 wherein the disk interface module comprises a  
2 data memory for temporarily storing a data file that is to be written to disk and a  
3 parity generator for generating parity information on the data file stored in the  
4 data memory prior to the transfer of the data file to the disk.
- 1 31. The storage system of claim 30 wherein data and its associated parity are stored  
2 on disk in their entirety in order to avoid reading previously stored data to  
3 determine data modifications and re-calculating parity.
- 1 32. The storage system of claim 12 wherein the host interface module stores data to  
2 disk without regard to where any earlier versions of that data were previously  
3 stored and without regard to the file and file system to which the data belongs.

- 1 33. The storage system of claim 12 wherein the data to be stored on the plurality of  
2 disks is arranged in a plurality of data pages and wherein the storage system  
3 comprises a reference counter for maintaining reference counts on each data  
4 page so that unused disk space can be readily identified.
- 1 34. The storage system of claim 12 in which the modules are interconnected using a  
2 high-speed, non-blocking, crossbar switch for transferring data.
- 1 35. The storage system of claim 12 in which separate, serial Interprocessor  
2 Communication (IPC) channels are used to transfer metadata between pairs of  
3 modules, thereby enabling the crossbar switch to be used at maximum efficiency  
4 for transferring data and allowing memory elements in the external interface to be  
5 partitioned into dedicated incoming and outgoing data buffers.
- 1 36. A fault-tolerant computer system for providing scaleable data processing and  
2 storage services to a client from a plurality of storage resources, comprising  
3 a plurality of identical resource interface modules connected to the storage  
4 resources;  
5 a plurality of identical access interface modules which receive requests for  
6 service from the client and select a subset of the plurality of resource interface  
7 modules to provide the requested service and distribute the workload across the  
8 plurality of storage resources; and  
9 a switch for temporarily connecting one of the plurality of access interface  
10 modules to the selected subset of the plurality of resource interface modules for  
11 providing the service to the client, the switch being constructed in two identical  
12 halves which are interconnected so that a failure in one switch half does not  
13 make the computer system inoperative.
- 1 37. The computer system of claim 36 wherein one switch half is designated as active  
2 and the other switch half is designated as standby and the active switch half is

3 used to temporarily connect one of the plurality of access interface modules to  
4 the selected subset of the plurality of resource interface modules.

1 38. The computer system of claim 37 wherein the active switch half and the standby  
2 switch half exchange roles if the active switch half fails.

1 39. The computer system of claim 36 wherein each storage resource is connected to  
2 at least two resource interface modules so that a failure in any resource interface  
3 module does not prevent access to the each storage resource.

1 40. The computer system of claim 36 wherein each access interface module can  
2 assume the workload of any of the plurality of access interface modules so that a  
3 failure in any access interface module can be bypassed by assigning the  
4 workload of the failed module to another of the access interface modules.

1 41. The computer system of claim 36 wherein each resource interface module can  
2 assume the workload of any of the plurality of resource interface modules so that  
3 a failure in any resource interface module can be bypassed by assigning the  
4 workload of the failed module to another of the resource interface modules.

1 42. The computer system of claim 36 wherein data is stored on one of the plurality of  
2 storage resources and the computer system further comprises a parity generator  
3 which computes a data tag, including parity information, from the data and  
4 stores the data tag on the plurality of storage resource apart from the data.

1 43. The computer system of claim 36 wherein data is stored on one of the plurality of  
2 storage resources and an acknowledgement is returned to the client after the  
3 data has been stored and wherein the data is stored on at least two separate  
4 storage resources before the acknowledgement is returned to the client.

1 44. A method for providing high-performance, scalable data processing and storage  
2 services to a client from a plurality of resources, the method comprising  
3 (a) providing an access interface module which receives requests for service  
4 from the client;  
5 (b) using the access interface module to select a subset of the plurality of  
6 resources to provide the requested service and distribute the workload  
7 across the plurality of resources; and  
8 (c) using a switch fabric to temporarily connect the access interface module to  
9 the selected subset of the plurality of resources for providing the service to  
10 the client.

1 45. The method of claim 44 wherein step (b) comprises selecting the subset of the  
2 plurality of resources based on the relative demand placed on the subset of  
3 resources.  
4

1 46. The method of claim 44 wherein step (c) comprises:  
2 (c1) using a control switch fabric for transferring control information; and  
3 (c2) using a separate data switch fabric for transferring data.  
4

1 47. The method of claim 46 wherein step (c1) comprises optimizing the control  
2 switch fabric for transferring control information and step (c2) comprises  
3 optimizing the data switch fabric for transferring data.  
4

1 48. The method of claim 46 wherein the request for service includes control  
2 information and data and wherein step (b) comprises separating the control  
3 information and the data and step (c) comprises transferring the data to the  
4 selected subset of resources over the data switch fabric.



- 1 49. The method of claim 46 wherein step (c1) comprises using a non-blocking  
2 crossbar switch for data transfer and step (c2) comprises using an Ethernet  
3 switch for control information transfer.
- 1 50. The method of claim 44 further comprising:  
2 (d) providing a resource module connected to the plurality of resources; and  
3 (e) using the resource module to generate preallocation information that  
4 preallocates services from the plurality of resources in order to evenly  
5 distribute a workload across the plurality of resources.
- 1 51. The method of claim 50 wherein step (c) comprises connecting the access  
2 interface module to the resource module so that the resource module can  
3 transfer the preallocation information to the access interface module.
- 1 52. The method of claim 51 wherein step (b) comprises selecting a subset of the  
2 plurality of resources based on the preallocation information.
- 1 53. The method of claim 44 wherein step (b) comprises temporarily storing  
2 information transferred between the access interface module and the selected  
3 subset of the plurality of resources.
- 1 54. The method of claim 44 wherein step (a) further comprises providing a plurality of  
2 access interface modules each access interface module receiving service  
3 requests from a plurality of clients.
- 1 55. A method for providing high-performance, scalable storage services to a client  
2 from a plurality of disks, comprising  
3 (a) providing a disk interface module connected to the plurality of disks for  
4 controlling data stored on the plurality of disks;

- 5 (b) providing a host interface module which receives requests for storage  
6 service from the client and selects a subset of the plurality of disks to  
7 provide the requested storage and distribute the workload across the  
8 plurality of disks; and  
9 (c) using a switch fabric to temporarily connect the host interface module to  
10 the resource module for providing the storage service to the client.

- 1 56. The method of claim 55 wherein step (c) comprises:  
2 (c1) using a control switch fabric optimized for transferring control information;  
3 and  
4 (c2) using a separate data switch fabric optimized for transferring data.

- 1 57. The method of claim 56 wherein the request for service includes control  
2 information and data and wherein step (b) comprises separating the control  
3 information and the data and step (c) comprises transferring the data to the  
4 selected subset of resources over the data switch fabric.

- 1 58. The method of claim 56 wherein step (c1) comprises using a non-blocking  
2 crossbar switch for data transfer and step (c2) comprises using an Ethernet  
3 switch for control information transfer.

- 1 59. The method of claim 55 wherein step (a) comprises generating preallocation  
2 information that preallocates physical storage in the plurality of disks in order to  
3 evenly distribute data across the plurality of disks.

- 1 60. The method of claim 59 wherein the physical storage in the plurality of disks is  
2 divided into zones and step (a) further comprises preallocating selected zones to  
3 frequently-accessed data, wherein the selected zones are selected in order to  
4 decrease disk access time.

- 1 61. The method of claim 59 wherein step (b) comprises logically mapping data items  
2 to be stored into allocation units preallocated to the host interface module by the  
3 disk interface modules.
- 1 62. The method of claim 59 wherein the host interface module comprises a first data  
2 memory and the disk interface module comprises a second data memory and  
3 wherein step (a) comprises using the first data memory to temporarily store  
4 information transferred from the host interface module to the disk interface  
5 module and step (b) comprises using the second data memory to temporarily  
6 store information received by the disk interface module from the host interface  
7 module.
- 1 63. The method of claim 55 wherein step (b) further comprises providing a plurality of  
2 host interface modules, each host interface module receiving service requests  
3 from a plurality of clients.
- 1 64. The method of claim 55 wherein step (a) comprises providing a plurality of disk  
2 interface modules connected to the plurality of disks.
- 1 65. The method of claim 64 wherein the host interface module directs writes to disk  
2 interface modules of data based on the relative demand placed on those  
3 modules without regard to either the contents of associated disks or the location  
4 of other copies of the data.
- 1 66. The method of claim 65 further comprising  
2 (d) providing a metadata module that is connected to the switch fabric and  
3 assigns each data object an identifying handle.
- 1 67. The method of claim 66 further comprising:

2 (f) using the host interface module and the metadata module to store  
3 metadata information associated with each data object in the plurality of  
4 disks separately from, and independently of, data associated with each  
5 data object.

1 68. The method of claim 66 wherein the handle contains the location of the metadata  
2 in the plurality of disks.

1 69. The method of claim 66 wherein the metadata module is dedicated exclusively to  
2 fetching, caching and manipulating file and directory attributes of each data  
3 object thereby allowing data paths to be optimized to maximize data flow.

1 70. The method of claim 55 wherein step (a) comprises using the disk interface  
2 module to preallocate storage in the plurality of disks in allocation units.

1 71. The method of claim 70 wherein step (a) comprises using the host interface  
2 module to select a subset of the plurality of disks from the preallocated storage.

1 72. The method of claim 70 wherein step (a) comprises using the host interface  
2 module to read a file from the plurality of disks and write the file to disk in a single  
3 allocation unit.

1 73. The method of claim 70 wherein step (a) comprises using a data memory for  
2 temporarily storing a data file that is to be written to disk and using a parity  
3 generator to generate parity information on the data file stored in the data  
4 memory prior to the transfer of the data file to the disk.

1 74. The method of claim 73 wherein data and its associated parity are stored on disk  
2 in their entirety in order to avoid reading previously stored data to determine data  
3 modifications and re-calculating parity.

1 75. The method of claim 55 wherein step (b) comprises using the host interface  
2 module to store data to disk without regard to where any earlier versions of that  
3 data were previously stored and without regard to the file and file system to  
4 which the data belongs.

1 76. The method of claim 55 wherein the data to be stored on the plurality of disks is  
2 arranged in a plurality of data pages and wherein the method further comprises:  
3 (g) maintaining reference counts on each data page so that unused disk  
4 space can be readily identified.

1 77. The method of claim 55 wherein step (c) comprises interconnecting the modules  
2 using a high-speed, non-blocking, crossbar switch for transferring data.

1 78. The method of claim 55 wherein step (c) further comprises using separate, serial  
2 Interprocessor Communication (IPC) channels to transfer metadata between  
3 pairs of modules, thereby enabling the crossbar switch to be used at maximum  
4 efficiency for transferring data and allowing memory elements in the external  
5 interface to be partitioned into dedicated incoming and outgoing data buffers.

1 79. A method for providing fault-tolerant scaleable data processing and storage  
2 services to a client from a plurality of storage resources, comprising  
3 (a) providing a plurality of identical resource interface modules connected to  
4 the storage resources;  
5 (b) providing a plurality of identical access interface modules which receive  
6 requests for service from the client and select a subset of the plurality of  
7 resource interface modules to provide the requested service and distribute  
8 the workload across the plurality of storage resources; and  
9 (c) using a switch for temporarily connecting one of the plurality of access  
10 interface modules to the selected subset of the plurality of resource

11 interface modules for providing the service to the client, the switch being  
12 constructed in two identical halves which are interconnected so that a  
13 failure in one switch half does not make the computer system inoperative.

1 80. The method of claim 79 wherein step (c) comprises designating one switch half  
2 as active and designating the other switch half as standby and using the active  
3 switch half to temporarily connect one of the plurality of access interface modules  
4 to the selected subset of the plurality of resource interface modules.

1 81. The method of claim 79 further comprising:  
2 (d) exchanging the roles of the active switch half and the standby switch half if  
3 the active switch half fails.

1 82. The method of claim 79 wherein step (a) comprises connecting each storage  
2 resource to at least two resource interface modules so that a failure in any  
3 resource interface module does not prevent access to the each storage resource.

1 83. The method of claim 79 wherein step (b) comprises providing access interface  
2 modules that can assume the workload of any other of the plurality of access  
3 interface modules so that a failure in any access interface module can be  
4 bypassed by assigning the workload of the failed module to another of the  
5 access interface modules.

1 84. The method of claim 79 wherein step (a) comprises providing resource interface  
2 modules that can assume the workload of any other of the plurality of resource  
3 interface modules so that a failure in any resource interface module can be  
4 bypassed by assigning the workload of the failed module to another of the  
5 resource interface modules.

1 85. The method of claim 79 wherein data is stored on one of the plurality of storage  
2 resources and the method further comprises:  
3 (e) computing a data tag, including parity information, from the data; and  
4 (f) storing the data tag on the plurality of storage resource apart from the  
5 data.

1 86. The method of claim 79 wherein data is stored on one of the plurality of storage  
2 resources and an acknowledgement is returned to the client after the data has  
3 been stored and wherein the method further comprises:  
4 (g) storing data on at least two separate storage resources before the  
5 acknowledgement is returned to the client.

1 87. The method of claim 79 wherein data is stored on one of the plurality of storage  
2 resources and the method further comprises:  
3 (h) computing a data tag, including parity information, from the data; and  
4 (i) storing the data tag on the plurality of storage resource apart from the  
5 data.